

MAINTENANCE METHOD OF RETRIEVAL SYSTEM AND RETREIVAL SYSTEM

Patent number: JP2001337973 (A)

Publication date: 2001-12-07

Inventor(s): MATSUI KAZUKI; IWAYAMA NOBORU; KIHARA HIDETO +

Applicant(s): FUJITSU LTD +

Classification:

- international: **G06F12/00; G06F13/00; G06F17/30; G06F12/00; G06F13/00; G06F17/30; (IPC1-7): G06F12/00; G06F13/00; G06F17/30**

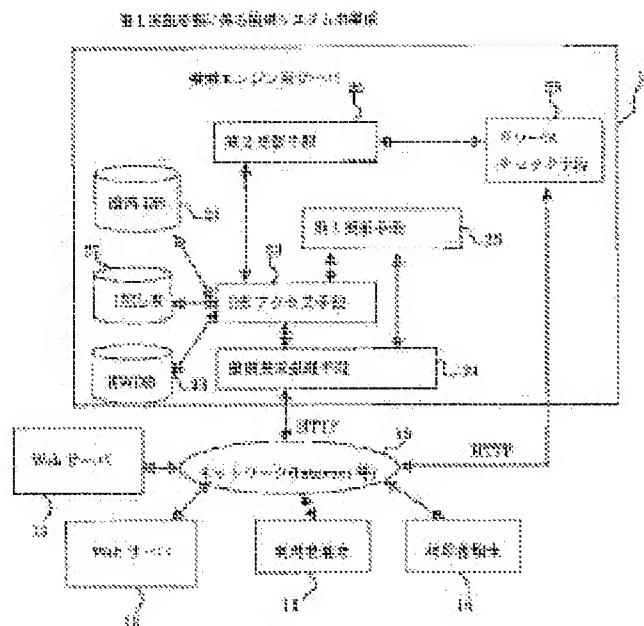
- european:

Application number: JP20000154502 20000525

Priority number(s): JP20000154502 20000525

Abstract of JP 2001337973 (A)

PROBLEM TO BE SOLVED: To enhance the retrieval precision in a retrieval system. **SOLUTION:** A retrieval DB is retrieved by use of a keyword ranked higher in a KWDB 22 to check the presence of link disconnection for the web page of the retrieval result. When a link disconnection occurs as the result of the check, the corresponding URL is eliminated from the retrieval DB 21. This processing is performed more frequently than other web pages. For a web page of the field conceivably attracting high users' attention, the link disconnection can be highly frequently checked, and the useless URL can be eliminated from the retrieval DB. Consequently, the retrieval precision can be efficiently enhanced.



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号
特開2001-337973
(P2001-337973A)

(43) 公開日 平成13年12月7日 (2001.12.7)

(51) Int.Cl. ⁷	識別記号	F I	テーマコード* (参考)
G 0 6 F 17/30	2 4 0	G 0 6 F 17/30	2 4 0 A 5 B 0 7 i
	1 1 0		1 1 0 F 5 B 0 8 2
	1 3 0		1 3 0 B
	3 4 0		3 4 0 B
12/00	5 4 6	12/00	5 4 6 B

審査請求 未請求 請求項の数 8 O L (全 13 頁) 最終頁に続く

(21) 出願番号 特願2000-154502(P2000-154502)

(22) 出願日 平成12年5月25日 (2000.5.25)

(71) 出願人 000003273

富士通株式会社

神奈川県川崎市中原区上小田中4丁目1番
1号

(72) 発明者 松井 一樹

神奈川県川崎市中原区上小田中4丁目1番
1号 富士通株式会社内

(72) 発明者 岩山 登

神奈川県川崎市中原区上小田中4丁目1番
1号 富士通株式会社内

(74) 代理人 100094145

弁理士 小野 由己男 (外2名)

最終頁に続く

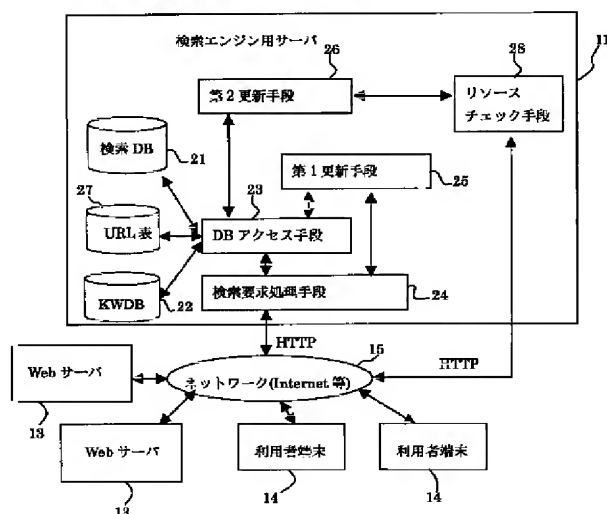
(54) 【発明の名称】 検索システムのメンテナンス方法及び検索システム

(57) 【要約】

【課題】 検索システムにおける検索精度を高める。

【解決手段】 KWDB 22において上位にランク付けされているキーワードを用いて検索DBを検索し、検索結果のウェブページについてリンク切れの有無をチェックする。チェックの結果、リンク切れが生じていれば、該当URLを検索DB 21から削除する。この処理を、他のウェブページよりも頻繁に行う。利用者に関心が高いと考えられる分野のウェブページについては、高い頻度でリンク切れをチェックでき、検索DBから無駄なURLを排除できる。その結果、効率よく検索精度を高めることが出来る。

第1実施形態に係る検索システムの構成



【特許請求の範囲】

【請求項1】ネットワーク上のリソースを利用者端末からの検索要求に応じて検索する検索システムのメンテナンス方法であって、

所定の基準で分類されたリソースを蓄積する検索蓄積手段と、リソースを検索するために用いられたキーワードを所定の順序でランク付けして蓄積するキーワード蓄積手段とを準備し、

リソースに順次アクセスしてリソースが無効になっているか否かをチェックし、かつキーワード蓄積手段で上位にランク付けされているキーワードにヒットするリソースについては前記チェックの頻度を高め、

リソースのチェック結果に基づいて検索蓄積手段を更新する、

検索システムのメンテナンス方法。

【請求項2】前記キーワード蓄積手段のキーワードは、検索に用いられた使用回数順にランク付けされている、請求項1に記載の検索システムのメンテナンス方法。。

【請求項3】前記キーワード蓄積手段のキーワードは、所定期間内におけるキーワードの使用回数の変化に基づいてランク付けされている、請求項1に記載の検索システムのメンテナンス方法。

【請求項4】禁止されている用語を蓄積した禁止用語蓄積手段をさらに準備し、

キーワード蓄積手段で上位にランク付けされているキーワードが禁止用語蓄積手段に蓄積されている用語に含まれる場合、当該キーワードについての前記チェックを行わない、請求項1～3のいずれかに記載の検索システムのメンテナンス方法。

【請求項5】所定の類似用語を同一の用語と見なしてキーワード蓄積手段を構築する、請求項1～4のいずれかに記載の検索システムのメンテナンス方法。

【請求項6】前記キーワード蓄積手段は、所定の分類基準に基づいて複数準備されている、請求項1～5のいずれかに記載の検索システムのメンテナンス方法。

【請求項7】ネットワーク上のリソースを利用者端末からの検索要求に応じて検索する検索システムであって、所定の基準で分類されたリソースを蓄積する検索蓄積手段と、

検索に用いられたキーワードを所定の順序でランク付けして蓄積するキーワード蓄積手段と、

利用者端末からの検索要求に基づいて検索蓄積手段を検索し、検索結果を要求元に通知する要求処理手段と、

検索毎にキーワード蓄積手段を更新する第1更新手段と、

リソースにアクセスしてリソースが無効になっているか否かをチェックし、かつキーワード蓄積手段で上位にランク付けされているキーワードにヒットするリソースについては、前記チェックを優先的に行うリソースチェック手段と、

リソースのチェック結果に基づいて、検索蓄積手段を更新する第2更新手段と、

を備えた検索システム。

【請求項8】ネットワーク上のリソースを利用者端末からの検索要求に応じて検索する検索システムであって、所定の基準で分類されたリソースを蓄積する検索蓄積手段と、

所定の分類基準ごとに設けられ、検索に用いられたキーワードを所定の順序でランク付けして蓄積する複数のキーワード蓄積手段と、

利用者端末からの検索要求に基づいて検索蓄積手段を検索し、検索結果を要求元に通知する要求処理手段と、

検索要求に基づいて更新すべきキーワード蓄積手段を特定し、特定したキーワード蓄積手段を検索毎に更新する更新手段と、

を備えた検索システム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、インターネットなどで利用できるWWW(World Wide Web)の検索システムに関する。

【0002】

【従来の技術】インターネット上には、各種の検索システムが存在している。これらの検索システムでは、キーワードなど所定の基準に基づいてウェブページを分類し、分類結果としてウェブページの存在場所を示すアドレス、例えばURL(Uniform Resource Locator)を検索DB(データベース)に蓄積している。そして、この検索DBに基づいて、入力されたキーワードにヒットするウェブページを検索する。

【0003】表1は、検索DBに蓄積される情報の概念を示す。一般的に、検索DBへの蓄積の仕方として、ウェブページのURLと、そのウェブページに記述されている文章等のコンテンツに含まれる単語等のキーワード情報とを対応付けた表を作る。そして、検索キーワードとともに検索要求が送られてくると、この表を検索してキーワードにマッチするURLを返す処理を行う。

【0004】

【表1】

「キーワードに基づいて分類されたURL表」の例

URL	キーワード
http://www.fujitsu.co.jp	富士通
http://www.flab.fujitsu.co.jp	富士通研究所
...	...

インターネット上の検索システムでは、通常このテーブルをインターネット上でアクセス可能な全てのウェブページについて作成する。よって、このテーブルは非常に膨大なデータ量に達するため、その維持・管理に多くの処理時間がかかる。

【0005】具体的には、ウェブページは、他のURLに移動したり、消滅したりして変化するので、検索システムではその変化を検索DBに反映させなければならない。そのため検索システムは、常にインターネット上のウェブページを巡回し、未登録のものを検索DBに登録すると共に、その変化を検出して検索DBを更新している。

【0006】検索DBの更新方法には、大別して2つの方法がある。1つは、メタ検索と呼ばれる方法である。この方法では、ある一つの検索エンジンが利用者の検索要求を受けると、複数の検索エンジンがその検索を実行する。そして、戻ってきた検索結果の中からリンク切れになっているものを探しだし、検索DBから削除する方法である（MetaCrawler等）。

【0007】他の1つの方法は、クライアント端末側でフィルタリング処理する方法である。すなわち、通常利用者の端末にインストールされている検索用ソフトウェアにより、検索結果の中でリンク切れのURLを削除してから表示するものである。

【0008】

【発明が解決しようとする課題】しかし、昨今のインターネットの利用拡大に伴い、ウェブページ数が飛躍的に増大し、検索システムがウェブページの変化についていけなくなっている。例えば、検索結果のURLにアクセスした結果、“Not Found（リンク切れ）”になったり、検索結果に同一のURLが何個もでたりする場合が多々ある。

【0009】従来の検索システムでは、元々の検索DB 21の中にリンク切れのウェブページの情報が残っている。そのため、検索速度が遅くなるとともに、利用者にとって不要なウェブページが検索結果に残ってしまい、検索の精度を下げることになっている。これは、検索DBの更新頻度が、ウェブページの更新頻度に比べて低いことに大きな原因がある。

【0010】なぜなら、リンク切れをチェックするためには、膨大な量のウェブページを一巡してそのウェブページの存在場所に変更がないかを確認しようとする、現在の検索システムの能力では2週間以上要することが通常である。一方、ウェブページは、毎日のように更新されるものが少なくない。その結果、検索DBには、実際のウェブページを反映していない情報が蓄積されたままとなり、検索結果の精度が下がってしまう。

【0011】一方、いくつかの検索システムでは、検索キーワード（KW）のランキング情報を蓄積したキーワードDBを公開している。キーワードDBのランキング

情報を見ると、一般的に検索キーワードは、検索に使用された時間が近い場合、類似する傾向にある。例えばその時に話題になっている製品の機種、ニュース、有名人の名前などが、検索キーワードのランキング上位になることが多い。しかし、検索システムにおいて、キーワードDBを活用して検索やその他の処理に用いた例はまだない。

【0012】本発明は、前記課題を解決するために、無駄な検索結果が含まれることを防止し、ウェブページの検索精度を向上させることを目的とする。また、本発明は、キーワードDBを活用して検索システム上で提供されるサービスの質をさらに向上させることを目的とする。

【0013】

【課題を解決するための手段】前記課題を解決するために、本願第1発明は、ネットワーク上のリソースを利用者端末からの検索要求に応じて検索する検索システムのメンテナンス方法であって、

A；所定の基準で分類されたリソースを蓄積する検索蓄積手段と、リソースを検索するために用いられたキーワードを所定の順序でランク付けして蓄積するキーワード蓄積手段とを準備し、

B；リソースに順次アクセスしてリソースが無効になっているか否かをチェックし、かつキーワード蓄積手段で上位にランク付けされているキーワードにヒットするリソースについては前記チェックの頻度を高め、

C；リソースのチェック結果に基づいて検索蓄積手段を更新する、

検索システムのメンテナンス方法を提供する。

【0014】ネットワーク上のリソースとは、例えばインターネット上のウェブページである。検索蓄積手段には、例えばキーワードに基づいて分類されたウェブページのURLが蓄積されている。キーワード蓄積手段には、ウェブページの検索に用いられたキーワードが、例えば使用回数順に蓄積されている。キーワード蓄積手段は、検索が行われるたびに更新される。

【0015】この検索システムは、ネットワーク上のウェブページに次々にアクセスし、そのウェブページが削除されていたり他に移動してリンク切れになっているかどうかをチェックする。このチェックをインターネット上の全てのウェブページについて行くと、各ウェブページについてのチェック頻度が例えば2～4週間に1回程度になってしまう。そこで、検索頻度の多いウェブページ、すなわちよく使用されるキーワードを含むウェブページについては、チェックを高い頻度、例えば1日1回行う。リンク切れが生じていると判定された場合、そのURLは無効なので検索蓄積手段から削除される。

【0016】インターネット上の全てのウェブページをもれなく同じ頻度でチェックするのではなく、使用頻度が高いと考えられるウェブページについてのみチェック

頻度を高くすることにより、ウェブページの検索速度と検索精度とを効率的に向上させる。本願第2発明は、前記第1発明において、前記キーワード蓄積手段のキーワードは、検索に用いられた使用回数順にランク付けされている検索システムのメンテナンス方法を提供する。

【0017】利用者により検索される頻度が高いウェブページについては、リンク切れを頻繁にチェックすることが好ましい。そこで、キーワードを使用頻度順にランク付けし、使用頻度が高いキーワードを含むウェブページについてチェック頻度を高める。本願第3発明は、前記第1発明において、前記キーワード蓄積手段のキーワードは、所定期間内におけるキーワードの使用回数の変化に基づいてランク付けされている検索システムのメンテナンス方法を提供する。

【0018】例えば、急激に使用回数が高まったキーワードは、利用者の関心の高まりを表していると考えられる。そのようなキーワードについては、総使用回数が多くなくても、そのキーワードを含むウェブページを頻繁にチェックすることが好ましいと考えられる。本願第4発明は、前記発明の方法において、禁止されている用語を蓄積した禁止用語蓄積手段をさらに準備し、キーワード蓄積手段で上位にランク付けされているキーワードが禁止用語蓄積手段に蓄積されている用語に含まれる場合、当該キーワードについての前記チェックを行わない検索システムのメンテナンス方法を提供する。

【0019】例えば、差別用語や放送禁止用語などを、禁止用語蓄積手段に蓄積しておく。このような用語については、例えキーワード蓄積手段で上位になっても、リンク切れのチェックを行わない。本願第5発明は、前記発明の方法において、所定の類似用語を同一の用語と見なしてキーワード蓄積手段を構築する検索システムのメンテナンス方法を提供する。

【0020】例えば既存の類似用語辞書を用いて類似用語群を同一キーワードとみなす。例えば、空白や長音の有無、平仮名やカタカナ、ローマ字の違いのみしかない用語を、同一のキーワードと見なし、使用回数を算出する。実情に即したキーワード蓄積手段の構築が可能となり、ひいては検索精度の向上を高めることが出来る。本願第6発明は、前記発明の方法において、前記キーワード蓄積手段が所定の分類基準に基づいて複数準備されている、検索システムのメンテナンス方法を提供する。例えば、ドメイン名や検索の時間帯ごとにキーワード蓄積手段を設けておく。

【0021】本願第7発明は、ネットワーク上のリソースを利用者端末からの検索要求に応じて検索する検索システムであって、検索蓄積手段と、キーワード蓄積手段と、要求処理手段と、第1更新手段と、リソースチェック手段と、第2更新手段と、を備えた検索システムを提供する。検索蓄積手段は、所定の基準で分類されたリソースを蓄積する。キーワード蓄積手段は、検索に用いら

れたキーワードを所定の順序でランク付けして蓄積する。要求処理手段は、利用者端末からの検索要求に基づいて検索蓄積手段を検索し、検索結果を要求元に通知する。第1更新手段は、検索毎にキーワード蓄積手段を更新する。リソースチェック手段は、リソースにアクセスしてリソースが無効になっているか否かをチェックし、かつキーワード蓄積手段で上位にランク付けされているキーワードにヒットするリソースについては、前記チェックを優先的に行う。第2更新手段は、リソースのチェック結果に基づいて、検索蓄積手段を更新する。

【0022】第1発明のメンテナンス方法を適用した検索システムである。本願第8発明は、ネットワーク上のリソースを利用者端末からの検索要求に応じて検索する検索システムであって、検索蓄積手段と、キーワード蓄積手段と、要求処理手段と、更新手段と、を備えた検索システムを提供する。検索蓄積手段は、所定の基準で分類されたリソースを蓄積する。キーワード蓄積手段は、所定の分類基準ごとに複数設けられ、検索に用いられたキーワードを所定の順序でランク付けして蓄積する。要求処理手段は、利用者端末からの検索要求に基づいて検索蓄積手段を検索し、検索結果を要求元に通知する。更新手段は、検索要求に基づいて更新すべきキーワード蓄積手段を特定し、特定したキーワード蓄積手段を検索毎に更新する。

【0023】この検索システムでは、キーワード蓄積手段が、所定の分類ごとに複数設けられている。例えば、ドメイン名、利用者の年代や居住地域、利用時間帯などにより、キーワード蓄積手段を分類する。例えばドメイン名毎にキーワード蓄積手段を設けた場合、各キーワード蓄積手段は、該当するドメイン名を有する利用者端末から検索要求があったときに更新される。

【0024】

【発明の実施の形態】＜発明の概要＞本発明は、インターネット上のウェブページ（リソース）を検索する検索システムにおいて、検索DBを効率よく更新する。具体的には、検索される頻度が高いウェブページについては、他のウェブページよりもリンク切れのチェックを頻繁に行う。そのために、検索システムで構築されるキーワードDBを活用する。

【0025】例えば、キーワードDBで使用回数が上位にランク付けされるキーワードを用いて検索DBを検索した場合、検索結果としてURLの集合（以下、「キーワードにヒットする」という）が得られる。これらのURLは、頻繁に検索されるウェブページのURLと考えられる。そこで、このようなキーワードにヒットするウェブページについては、リンク切れのチェックを例えば1日1回行う。チェックの結果、リンク切れとなっているURLは検索DBから削除される。

【0026】上位のキーワードにヒットしないウェブページについては、従来通り例えば2週間に1回程度のチ

ェック頻度のままであるものの、ユーザにとってチェックの必要性が高いウェブページについて優先的にリンク切れをチェックするので、効率的に検索DBを更新することが出来る。結果として、検索DBから無効なURLが削除され、検索の速度及び精度が向上する。

【0027】＜第1実施形態例＞

〔構成〕次に、本発明の検索システムについて、具体的に説明する。図1は、第1実施形態例に係る検索システムの全体構成図である。検索システムは、検索エンジン用サーバ（以下、単に検索サーバという）11、ウェブサーバ13及び利用者端末14が、インターネット15で接続されて構成されている。

【0028】ウェブサーバ13は、ウェブページを収集し、検索サーバ11や利用者端末14に提供している。利用者端末14には、ブラウザ機能が搭載され、ウェブサーバ13が提供するウェブページを表示可能である。また、利用者端末14には、検索サーバ11に検索要求を送信し、検索結果を出力するための検索クライアントが搭載されている。

【0029】検索サーバ11は、検索DB21、キーワードDB（KWDB）22、DBアクセス手段23、検索要求処理手段24、第1更新管理手段25、第2更新手段26、URL表27及びリソースチェック手段28を有している。検索DB21には、ウェブページのURLが、キーワードに基づいて分類されて蓄積されている。検索DB21には、希望者がURLを登録する登録型と、自動的にURLを収集して情報を構築するソフトウェアロボット収集型とがあり、そのどちらにも本発明を適用可能である。

【0030】KWDB22には、所定の順序でランク付けられたキーワードが蓄積されている。本実施形態例では、キーワードは、検索に使用された使用回数順にランク付けされている。図2は、KWDB22の概念説明図である。KWDB22には、検索キーワードと、そのキーワードが検索に使用された使用回数と、順位とが、記憶されている。本発明では、このKWDB22を利用してKWDB22の中で例えば上位10位までについて、リンク切れなどのチェックを他のキーワードよりも短い時間間隔で行うようにする。例えば、通常は週1回の頻度でしかチェックされない場合には、上位のキーワードについては毎日1回チェックを行うようにする。

【0031】検索要求処理手段24は、利用者端末14からの検索要求を受信し、キーワードにヒットするウェブページのURLの集合を、検索結果として利用者端末14に送信する。DBアクセス手段23は、利用者端末14からの検索要求に従って検索DB21を検索する。また、DBアクセス手段23は、第1更新手段25からの指示に従い、KWDB22を更新する。さらにDBアクセス手段23は、第2更新手段26からの指示に従い、リンク切れしたURLを、検索DB21から削除す

る。

【0032】第1更新手段25は、検索DB21の検索に用いられたキーワードに基づいて、KWDB22内のキーワードの順序を並べ替える。例えば、第1更新手段25は、KWDB22内のキーワードの新たな順位を計算し、キーワードの並べ替えをDBアクセス手段23に指示する。第2更新手段26は、リソースチェック手段28を用いてインターネット上のウェブページを巡回し、リンク切れを検出する。また、第2更新手段26は、KWDB22で上位にあるキーワードを用いて検索DB21を検索し、検索結果にヒットしたURLについてリソースチェック手段28を用いてリンク切れしていないかどうかのチェックを行う。このチェックは、他のウェブページにおけるリンク切れの検出よりも高い頻度で行われる。また、第2更新手段26は、重複するURLが検索結果に出現しないか否かのチェックも行ふ。これらのチェック結果に基づいて、第2更新手段26は、検索DB21の更新を行う。

【0033】図3は、URL表27の概念説明図である。検索DB21には、キーワードに基づいて分類されたURLが蓄積されているが、その他にURL表を設けておくとうまい。URL表(a)では、チェック時間、検索キーワード、URL表(b)へのポインタが記憶されている。「チェック時間」は、対応する検索キーワードにヒットするURLについてチェックした時間である。所定時間以上経過したエントリは、このチェック時間に基づいて削除される。「検索キーワード」は、検索に用いられたキーワードである。「URL表(a)へのポインタ」は、対応するキーワードで検索した場合にヒットするURLのリンク切れ検出回数が格納されているアドレスを示す。

【0034】URL表(b)では、リンク切れ検出回数と、URLとが記憶されている。「リンク切れ検出回数」は、対応するURLについてリンク切れが検出された回数を示している。この例は、“HYPERLINK http://www.fujitsu.co.jp www.fujitsu.co.jp”についてはリンク切れが検出されていないが、“www.fujitsu.co.jp”については3回のリンク切れが検出されていることを示している。

【0035】リソースチェック手段28は、検索サーバ11から指示されたウェブページのデータを収集し、リンク切れが生じているか否かをチェックする。また、リソースチェック手段28は、チェック結果を第2更新手段26に通知する。

〔処理方法〕

(1) 検出処理

図4は、検索システムが行う検出処理の流れの一例を示すフローチャートである。この処理では、検索システムは、KWDB22で所定の順位のキーワードについて、リンク切れ及び重複URLがあるか否かを、所定時間ご

とにチェックする。説明を容易にするために、1位のキーワードについて、1日1回チェックする場合を例に取る。

【0036】ステップS1；第2更新手段26は、KWDB22で1位にランク付けされているキーワードを、DBアクセス手段23を介して取得する。

ステップS2；第2更新手段26は、取得したキーワードを用いて検索DB21の検索を行い、このキーワードにヒットするURLの集合を取得する。

ステップS3；第2更新手段26は、取得したURLの集合の中から、重複するURLを削除する。

【0037】ステップS4；第2更新手段26は、取得した検索結果に基づいて、URL表を更新する。すなわち、検索に用いたキーワード及び取得したURLのためのエントリを、URL表(a)及び(b)に設ける。既にURLがURL表に登録されている場合は、そのURLについては新たなエントリを設けない。

ステップS5、S6；第2更新手段26は、取得したURLについて、リンク切れが生じているか否かを、リソースチェック手段28を用いてチェックする(S5)。リンク切れが生じていると、そのURLのリンク切れ検出回数をインクリメントする(S6)。

【0038】ステップS7；第2更新手段26は、検索したキーワードにヒットする全てのURLについて、リンク切れをチェックしたか否かを判断する。“Yes”と判断すると、ステップS8に移行する。“No”と判断すると、再び前記ステップS5に戻り、次のURLについてリンク切れをチェックする。

ステップS8；第2更新手段26は、所定時間T1、本例では1日が経過するのを待機し、前記ステップS1に戻って前述の処理を繰り返す。

【0039】本例では、KWDB22で1位のキーワードについてのみ前述の処理を行ったが、例えば所定順位を100位とする場合、各キーワードについて順次前述の処理を行えばよい。すなわち、第2更新手段26は、1日に1回、KWDB22で1位のキーワードにヒットするURLについて、リンク切れ及び重複URLがあるかどうかをチェックし、その結果をURL表(a)

(b)に書き込む処理を行う。KWDB22で1位のキーワードにヒットしない他のウェブページについては、従来通り巡回して順次リンク切れをチェックする。後述するように、1位のキーワードについては、検索DB21からリンク切れURLが削除されるので、検索速度及び精度を高めることが出来る。

【0040】さらに、このURL表に基づいて検索結果を出力すれば、検索結果から重複URLが削除されているので、一層検索精度を向上することが出来る。この場合、第2更新手段26は、所定時間T1よりも長い時間間隔T3でURL表を参照し、チェック時間から所定時間が経過したエントリを削除する処理を行う。URL表は、

上位のキーワードによる検索結果を示しているものの、あまりにも古い検索結果を用いることを防止するためである。

【0041】(2)DB更新処理

図5は、第2更新手段26が行うDB更新処理の流れを示すフローチャートである。この処理では、所定時間T2($T2>T1$)毎に、リンク切れが所定回数以上のURLをURL表及び検索DB21から削除する処理を行う。

ステップS21；第2更新手段26は、所定時間T2が経過するのを待機し、経過するとステップS22に移行する。

【0042】ステップS22、S23、S24；第2更新手段26は、その時点で作成されているURL表

(b)を参照し、リンク切れ検出回数が所定数以上のURLを検索する。そのようなURLが存在する場合、URL表(a)(b)及び検索DB21から該当URLを削除する。

ステップS25；第2更新手段26は、URL表(b)の全てのエントリについて、リンク切れ検出回数をチェックしたか否かを判断す。“Yes”と判断すると前記ステップS21に戻って前述の処理を繰り返す。“No”と判断すると、URL表(b)の次のエントリのURLについて、リンク切れ検出回数をチェックする(S22)。

【0043】すなわち、この処理では、リンク切れ検出回数が所定数以上のURLを、URL表(a)(b)から削除し、またそのURLを検索DB21から削除する。これにより、無用なURLが検索DB21から削除されて検索処理速度が上がることを期待できる。さらに、無用なURLが検索結果に出現しなくなり、検索精度を高めることができる。

【0044】なお、リンク切れが1回でもあればURL表及び検索DB21から削除することも考えられる。しかし、諸原因、例えばネットワークのトラフィックやウェブサーバの一時的なトラブルによりリンク切れではないのにアクセスできない場合を考慮し、若干の余裕を持たせて削除を行うことが好ましい。

<第2実施形態例>図6は、第2実施形態例に係る検索システムの全体構成図である。この検索システムは、KWDB22が複数設けられている点を除き、前記第1実施形態例と同様の構成を有する。

【0045】この例では、複数のKWDB22、すなわちKWDB-a及びKWDB-bが、ドメイン名毎に設けられている。具体的には、KWDB-aは、利用者端末14のドメイン名が“co.jp”の場合に用いられる。KWDB-bは、利用者端末14のドメイン名が“ne.jp”の場合に用いられる。例えば、会社から仕事に関するウェブページを検索する場合と、自宅から趣味に関するウェブページを検索する場合とを考える。仕事に関するウェブページと趣味に関するウェブページとでは、内容やア

クセス頻度が異なる場合が多い。従って、同じKWDB 22を使用するのは適切ではない場合がある。

【0046】例えば、仕事に関する分野では、キーワードがKWDB 22の上位にランキングされ、URLのチェックが頻繁に行われる。一方、趣味に関する分野では、その分野ではよく使用されるキーワードであっても全体としてのキーワードの使用回数が少ないため、KWDB 22の上位にランク付けされない。結果として、この分野のウェブページについて、リンク切れチェックの頻度が低下する。

【0047】そこで、所定の分類、例えば利用者端末14のドメイン名毎にKWDB 22を設け、検索要求を行った利用者端末14のドメイン名に応じて使用するKWDB 22を切り替える。利用者端末14のドメイン名は、インターネットで用いられるhttpプロトコルに準拠して検索サーバ11に送信されるIPアドレスに通常含まれているので、そこから取得可能である。

【0048】図7は、ドメイン名毎にKWDB 22を更新する処理の流れを示すフローチャートである。この処理では、検索要求がある毎に、ドメイン名に応じてKWDB 22を更新する処理を行う。

ステップS31；検索要求処理手段24は、利用者端末14から検索要求を受信する。

【0049】ステップS32；検索要求処理手段24は、利用者端末14のIPアドレスなどからドメイン名を抽出する。

ステップS33；検索要求処理手段24は、利用者端末14のドメイン名に対応するKWDB 22を選択する。
ステップS34；検索要求処理手段24は、検索要求に含まれる検索キーワードを抽出する。

【0050】ステップS35；検索要求処理手段24は、DBアクセス手段23を用い、選択したKWDB 22に抽出した検索キーワードが含まれているか否かを判断する。含まれていれば、そのキーワードの使用回数をインクリメントとし、必要に応じて順位を変更する。含まれていなければ、そのキーワードをKWDB 22に登録し、使用回数を“1”とする。

【0051】ステップS36；検索要求処理手段24は、検索キーワードに基づいて検索DB 21の検索を実行し、検索結果を利用者端末14に送信する。

＜第3実施形態例＞前記第1及び第2実施形態例では、KWDB 22においてキーワードの使用回数順にキーワードをランク付けし、使用回数が多いキーワードについて優先的に検出処理及びDB更新処理を行っている。しかし、他の順序でキーワードをランク付けすることも考えられる。本実施形態例では、所定期間内における使用回数の変化に基づいて、優先的に検出処理及びDB更新処理を行う。

【0052】図8は、本実施形態例におけるKWDB 22の概念説明図である。このKWDB 22には、検索キ

ーワード、使用回数、順位、一定時間前の順位、変化速度、一定時間前の変化速度及び変化加速度が記録されている。このKWDB 22の順位は、一定時間、この例では一週間毎に再計算される。「順位」は、変化加速度の大きさに基づいて決定されている。「一定時間前の順位」とは、例えば一週間前の順位である。「変化速度」とは、「一定時間前の順位」から「順位」への変化を示す。「一定時間前の変化速度」とは、例えば一週間前に置ける変化速度である。「変化加速度」とは、「一定時間前の変化速度」から「変化速度」への変化率を示す。具体的には、先々週の変化速度から先週の変化速度への変化率を示す。

【0053】図8では、キーワード「富士通」の順位が1位になっている。このキーワードの一週間前の順位は2位である。先週の順位と今週の順位との差、すなわち「変化速度」は「1」となる。一週間前の時点における変化速度、すなわち「一定時間前の変化速度」は10となっている。先週と先々週との変化速度の差、すなわち「変化加速度」は「-9」となる。

【0054】このように、所定期間内における使用回数の変化に基づいてキーワードの順位を決定すると、例えば急激に関心が高まった分野のウェブページについてリンク切れなどのチェックをすることができる。従って、その時々ニーズに応じた検索DB 21の更新を行うことが可能となる。

＜第4実施形態例＞前記第3実施形態例においては、所定期間内における使用回数の変化に基づいて、KWDB 22内のキーワードの順位を決定している。本実施形態例では、使用回数の変化に加え、さらに前記図8の各項目に重み付けを行ってキーワードをランク付けする。

【0055】例えば、重み付けパラメータとして、「順位」に対し0.5、「変化速度」に対し0.3を用いる。これらの項目に重み付けパラメータを乗じた値に基づいて順位を算出する。具体的には、次のようにして算出する。図8のキーワード「富士通」の順位は1位であり、一定時間前の順位は2位である。先週からの順位の変化を示す変化速度は、 $2 - 1 * 0.5 = 1.5$ と計算される。また、変化加速度、すなわち先週と先々週との変化速度の差は、 $1.5 - 10 = -8.5$ と算出される。このようにして算出した変化加速度の値を用いてキーワードのランク付けを行う。

【0056】＜第5実施形態例＞第5実施形態例においては、放送禁止用語、差別用語などに相当する言葉にヒットするウェブページの変化を検索DB 21に反映させない。まず、放送禁止用語、差別用語などに相当する言葉を、禁止キーワードとして検索DB 21にあらかじめ登録しておく。検索サーバ11は、KWDB 22で上位のキーワードが禁止キーワードに該当するか否かの判断を行い、禁止キーワードに該当する場合、検索DB 21の更新処理を行わない。

【0057】図9は、検索サーバ11が行う処理の流れを示すフローチャートである。

ステップS41；検索要求処理手段24は、検索要求を利用者端末14から受信する。

ステップS42；検索要求処理手段24は、検索要求に含まれる検索キーワードを抽出し、第1更新手段25に通知する。

【0058】ステップS43；第1更新手段25は、抽出した検索キーワードが所定の禁止キーワードに含まれているか否かを判断する。“Yes”と判断すると後述するステップS45に移行し、そのまま通常の実行処理を実行する。“No”と判断すると、ステップS44に移行する。

ステップS44；第1更新手段25は、検索キーワードに基づいて、KWDB22を更新する。まず、抽出した検索キーワードがKWDB22に含まれているか否かを判断する。含まれていれば、そのキーワードの使用回数をインクリメントとし、必要に応じて順位を変更する。KWDB22に含まれていなければ、そのキーワードをKWDB22に新たに登録し、使用回数を“1”とする。

【0059】ステップS45；検索サーバ11は、検索キーワードに基づいて検索DB21を検索し、検索結果であるURLの集合を利用者端末14に送信する。この実施形態例を用いれば、公序良俗に反する内容のウェブページについてまで頻繁にリンク切れをチェックすることがなくなり、無用なチェックによる負担を省略でき、社会の道德に適合した検索DB21の更新を行うことが出来る。

【0060】＜第6実施形態例＞第6実施形態例では、KWDB22の構築にあたり、あらかじめ語句の類似関係を定義する類似用語辞書を検索DB21に登録しておく。KWDB22を更新する場合、類似範囲に属するキーワードは同一のキーワードとして扱う。図10は、類似用語辞書の概念説明図である。この図では、漢字、ローマ字、ひらがな、カタカナによる差しかない語句を、類似と見なす辞書の例を示している。類似用語辞書としては、通常用いられる辞書を使用できる。

【0061】本実施形態例に依れば、ひらがな、カタカナあるいはローマ字による違い、聴音や空白の有無の違いなどのみであり、実質的な違いのないキーワードをまとめて同一の語句としてKWDB22で扱うことが出来る。従って、KWDB22のランク付けが実情に即したものとなり、さらには検索DB21の更新が実情に即したものとなる。

【0062】＜第7実施形態例＞前記第2実施形態例においては、複数のKWDB22をドメイン名毎に構築している。しかし、他の分類に基づいて複数のKWDB22を構築することも可能である。例えば、個人情報を蓄積したユーザDBが検索システム上にあらかじめ存在し

ている場合、個人情報に基づいてKWDB22を構築しても良い。個人情報とは、例えばユーザ名、性別、居住地、年代、趣味などである。

【0063】例えば、年代に応じたKWDB22を構築しておき、検索要求があった場合、要求元ユーザの年代に対応するKWDB22を更新する。これにより、KWDB22のランク付けが各ユーザの年代に適合した内容となり、検索DB21は各年代のユーザの要求に合うように更新される。KWDB22を構築する他の分類として、アクセス時間毎にKWDB22を構築しておくことも考えられる。例えば、午前中はニュース関係、午後は株式関係、夜は趣味関係の検索が多いとする。その場合、午前中、午後、夜用のKWDB22を準備しておく。各KWDB22は、そのKWDB22の時間帯における検索要求に応じて更新される。

【0064】このようにすると、各KWDB22におけるキーワードの順位は、時間帯に応じた内容を反映したランク付けとなる。また、各自時間帯でアクセスの多い分野において、無効なURLを検索DB21から取り除くので、検索DB21を各時間帯における利用者のニーズに適合させることが出来る。

＜他の実施形態例＞

(A) 検索システムを複数のサーバで構成することも可能である。図11は、複数のサーバで構成された検索システムの構成例を示す。この図に示す検索システムは、前記検索サーバ11に相当する部分が、サーバ群Aと、リソースチェック手段を搭載したサーバ群Bと、検索DB21と、KWDB22と、URL表27とから構成されている。

【0065】サーバ群Aを構成するサーバには、DBアクセス手段23、検索要求処理手段24、第1更新手段25及び第2更新手段26が設けられている。また、サーバ群Bを構成するサーバには、リソースチェック手段28が設けられている。これらの機能は、前記第1実施形態例において述べたものと同様である。検索DB21、KWDB22及びURL表27には、前述の実施形態例で例示したような情報が蓄積される。これらは、それぞれ独立に存在してサーバ群A、Bから直接アクセスされたり、データベースサーバを介してサーバ群A、Bからアクセスされたりする。

【0066】(B) 前述した本発明の方法を実行するプログラムを記録した記録媒体は、本発明に含まれる。ここで記録媒体としては、コンピュータが読み書き可能なフロッピー（登録商標）ディスク、ハードディスク、半導体メモリ、CD-ROM、DVD、光磁気ディスク（MO）、その他のものが挙げられる。

【0067】

【発明の効果】本発明を用いれば、利用者の関心が高いと考えられる分野において優先的にリソースをチェックし、無駄になったリソースを検索蓄積手段から削除する

ことができる。従って、検索速度及び検索精度の向上を、利用者の関心が高い分野において期待することが出来る。

【図面の簡単な説明】

- 【図1】 第1実施形態例に係る検索システムの構成図。
- 【図2】 KWDBの概念説明図。
- 【図3】 URL表の概念説明図。
- 【図4】 検出処理の流れを示すフローチャート。
- 【図5】 DB更新処理の流れを示すフローチャート。
- 【図6】 第2実施形態例に係る検索システムの構成図。
- 【図7】 複数KWDBの更新処理の流れを示すフローチャート。

ャート。

【図8】 第3実施形態例におけるKWDBの概念説明図。

【図9】 禁止キーワードチェック処理の流れを示すフローチャート。

【図10】 類似キーワード辞書の概念説明図。

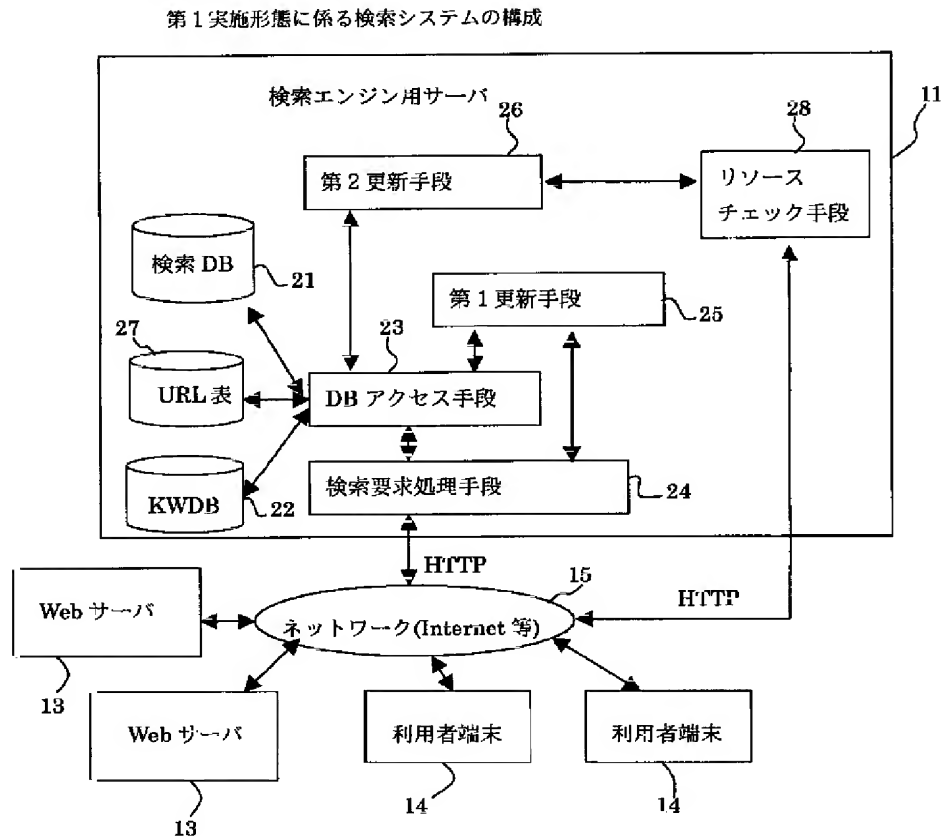
【図11】 複数のサーバで構成された検索システムの構成図。

【符号の説明】

21；検索DB

22；KWDB

【図1】



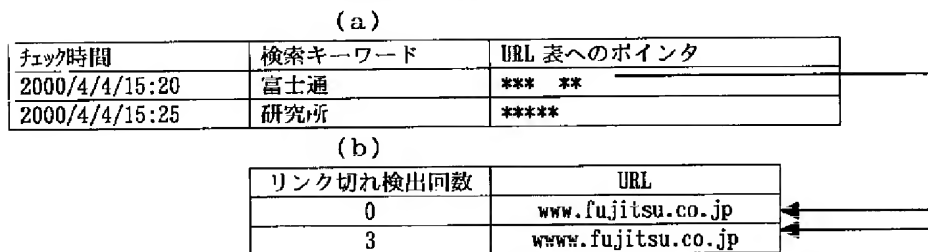
【図2】

KWDBの概念説明図

検索キーワード	使用回数	順位
富士通	123456	1
研究所	2345	?
...

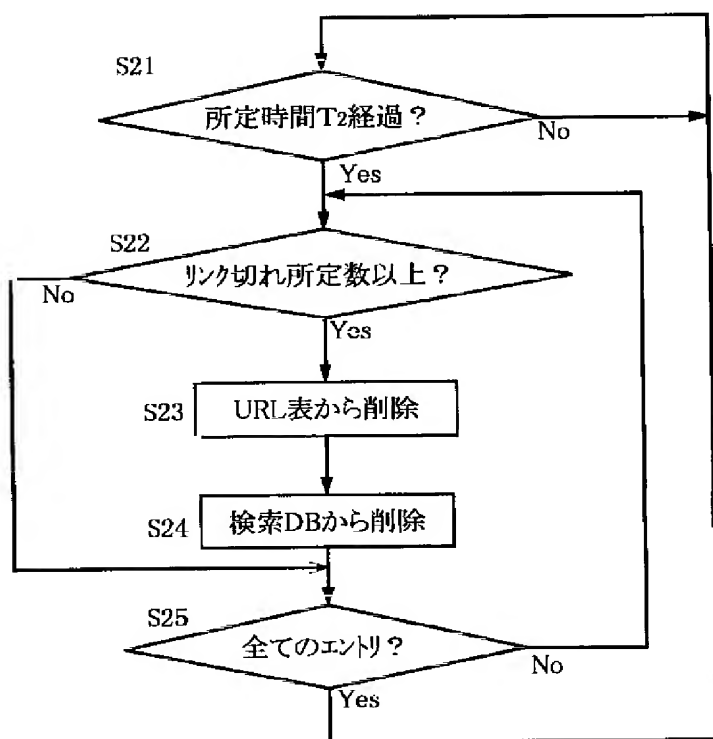
【図3】

URL表の概念説明図



【図5】

DB更新処理の流れ



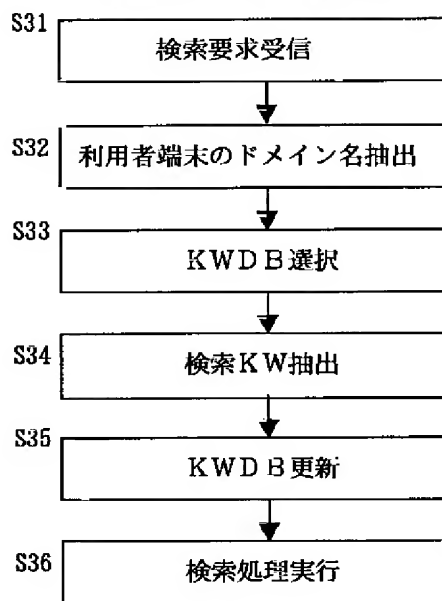
【図8】

第3実施形態例におけるKWDBの概念説明図

検索キーワード	使用回数	順位	一定時間前の順位	変化速度	一定時間前の変化速度	変化加速度
富士通	12345	1	2	1	10	-9
研究所	2345	2	12	10	5	5
...

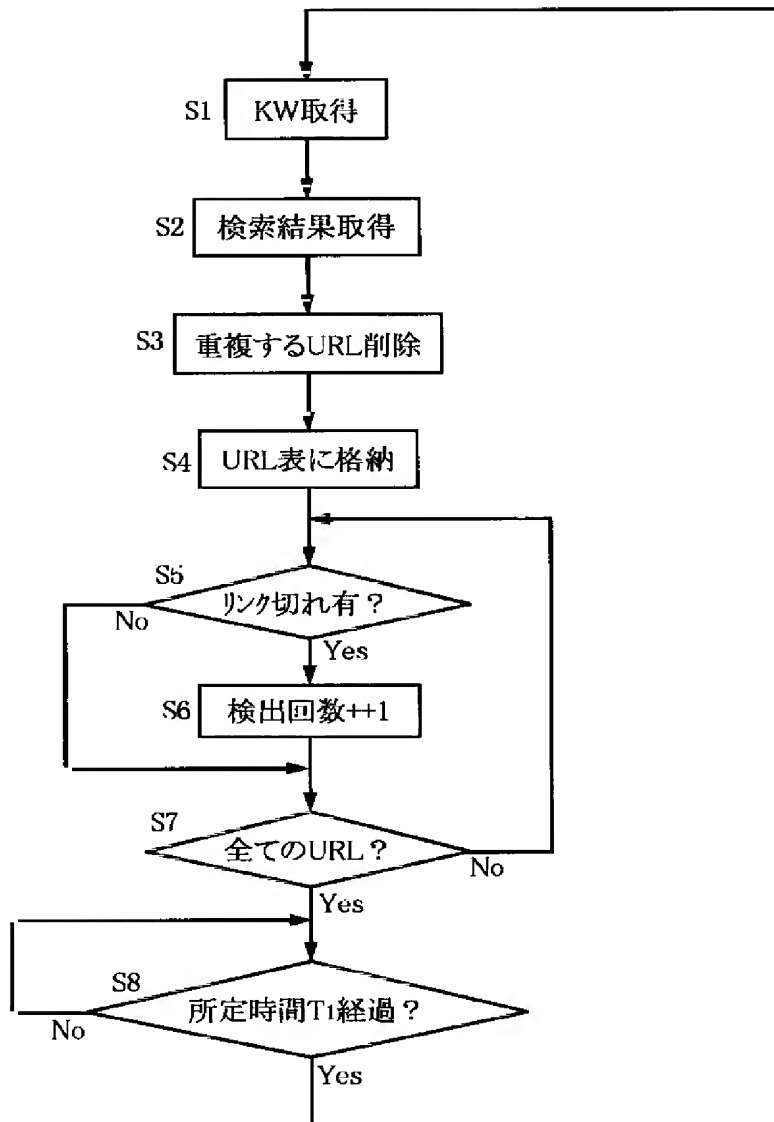
【図7】

複数KWDBの更新処理



【図4】

検出処理の流れ



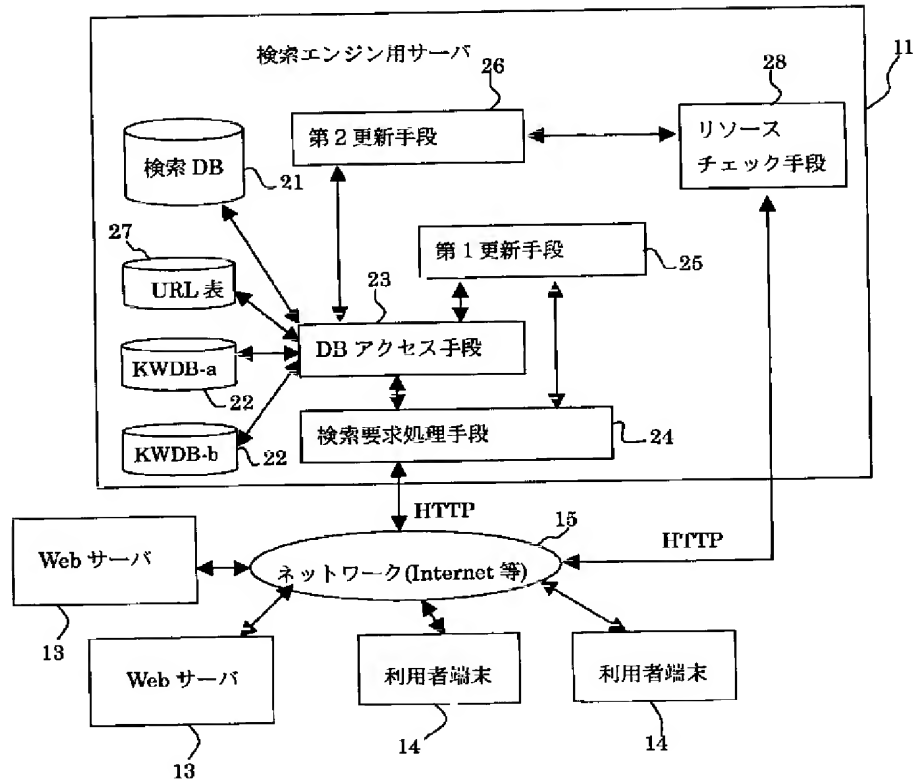
【図10】

類似キーワード辞書の概念説明図

検索キーワード	類似候補	類似候補	類似候補
富士通	FUJITSU	ふじつう	フジツウ
研究所	Laboratories	Laboratory	ラボ
...			...

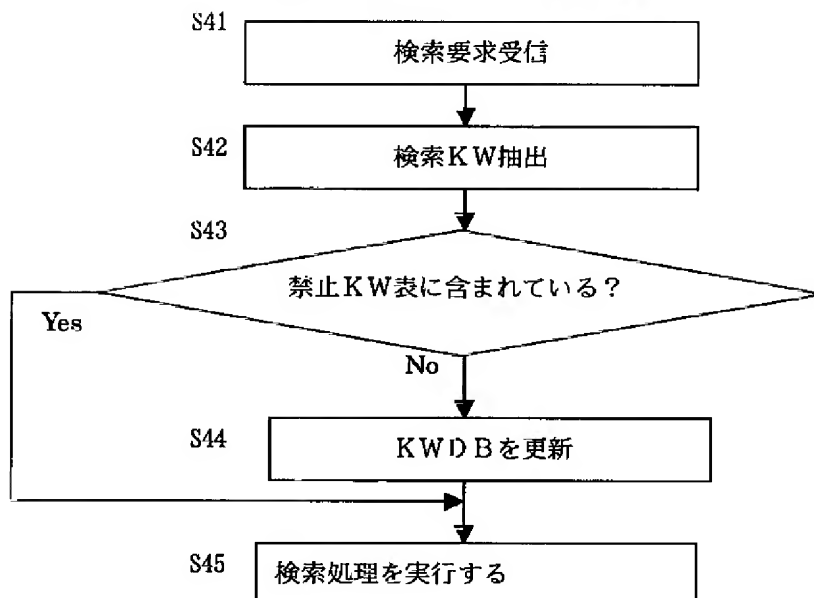
【図6】

第2実施形態に係る検索システムの構成

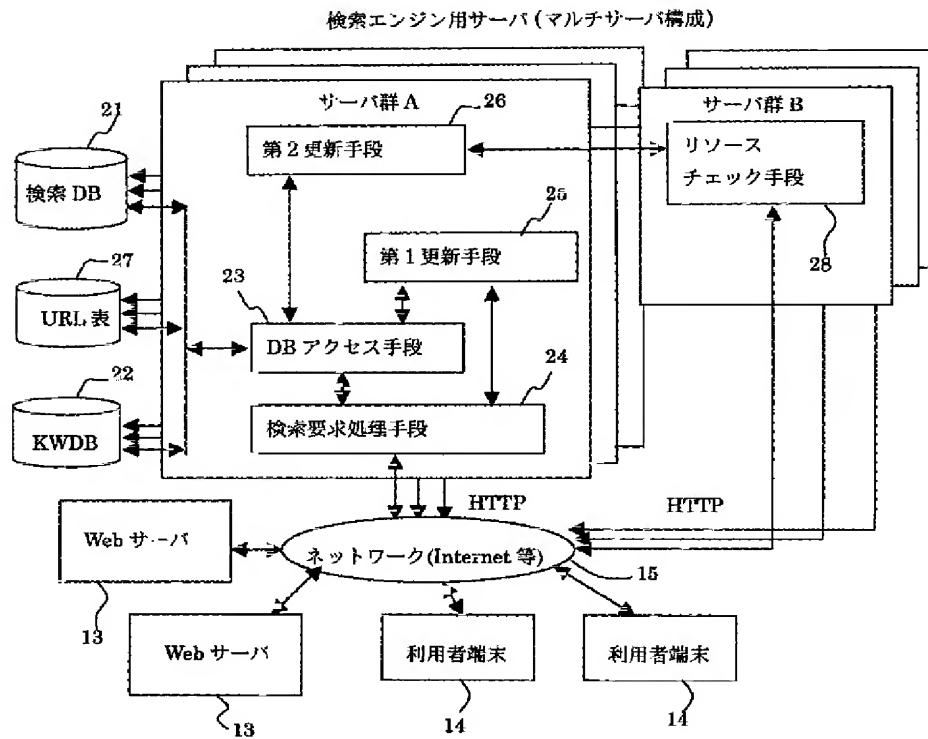


【図9】

禁止キーワードチェック処理の流れ



【図 1 1】



フロントページの続き

(51)Int.Cl. ⁷	識別記号	F I	(参考)
G 0 6 F 13/00	5 1 0	G 0 6 F 13/00	5 1 0 B
(72)発明者 木原 英人		F ターム(参考) 5B075 NK02 NR05 PR04	
神奈川県川崎市中原区上小田中4丁目1番		5B082 HA00	
1号 富士通株式会社内			